

Co-adaptation in a Handwriting Recognition System

Sunsern Cheamanunkul

University of California, San Diego
9500 Gilman Dr, La Jolla, CA 92093
scheaman@eng.ucsd.edu

Yoav Freund

University of California, San Diego
9500 Gilman Dr, La Jolla, CA 92093
yfreund@eng.ucsd.edu

ABSTRACT

Handwriting is a natural and versatile method for human-computer interaction, especially on small mobile devices such as smart phones. However, as handwriting varies significantly from person to person, it is difficult to design handwriting recognizers that perform well for all users. A natural solution is to use machine learning to adapt the recognizer to the user. One complicating factor is that, as the computer adapts to the user, the user also adapts to the computer and probably changes their handwriting. This paper investigates the dynamics of co-adaptation, a process in which both the computer and the user are adapting their behaviors in order to improve the speed and accuracy of the communication through handwriting. We devised an information-theoretic framework for quantifying the efficiency of a handwriting system where the system includes both the user and the computer. Using this framework, we analyzed data collected from an adaptive handwriting recognition system and characterized the impact of machine adaptation and of human adaptation. We found that both machine adaptation and human adaptation have significant impact on the input rate and must be considered together in order to improve the efficiency of the system as a whole.

Author Keywords

Co-adaptation; handwriting recognition; communication channel;

INTRODUCTION

Handwriting is a natural and versatile method for human-computer interaction, especially on small mobile devices such as smart phones. As handwriting varies significantly from person to person, it is difficult to design a handwriting recognition system that performs well for all users. Modern handwriting recognizers resort to machine learning techniques to adapt and specialize their handwriting models to each individual user. As the recognizer adapts to the human user, the user is likely to adapt to the system as well. We call this situation “co-adaptation” where both human and computer adapts to each other simultaneously.

In general, co-adaptation can manifest in any adaptive system. Designing a system that co-adapts with the users is a challenging problem on its own [1, 2, 3]. Our goal in this paper is not to address those challenges, but rather to focus on characterizing the impact of machine adaptation and of human adaptation in the context of handwriting recognition. We believe that this study will provide us with useful insights towards designing a more efficient adaptive handwriting recognition system.

In order to evaluate performance of a handwriting recognition system under co-adaptation, we introduce a framework based on the idea of Shannon’s communication channel [4] that considers both the user and the handwriting recognizer in a single system. Under this framework, we define the notion of “channel rate” that measures the amount of information successfully transferred from the user to the computer.

To quantify the effect of machine adaptation and of user adaptation empirically, we developed a handwriting recognition system that is capable of adapting to the handwriting of each individual user over time. We collected usage data from 15 different users and performed an analysis of the channel rate.

The paper is organized as follows. First, in Section , we present the information-theoretic framework for quantifying the efficiency of a handwriting system where the system includes both the user and the computer. Next, in Section , we describe our adaptive handwriting recognition algorithm that we developed for our experiment. Then, in Section , we describe the experiment and present the results in terms of the performance measures derived from the proposed framework. Finally, we draw some conclusions in Section .

HANDWRITING RECOGNITION AS A COMMUNICATION CHANNEL

Unlike typing, which transmits information to the computer at discrete time points, handwriting continuously transmits

information as the writer creates the trajectory. Traditionally, handwriting data is analyzed one “unit” at a time where “unit” can be a stroke, a character, a word or even a sentence. In this work, we propose an alternative analysis where the data is analyzed in fixed intervals of time. We consider the process of writing as a process through which the intended letter is disambiguated from the other possible letters.

We formalize this process using the concept of communication channel [4]. Let \mathcal{E} denote the set of all possible input. Technically, the set \mathcal{E} can be a set of sentences, a set of words, or a set of characters. Without loss of generality, in this work, we assume that \mathcal{E} is a set of 26 English characters. We also ignore dependencies between characters due to the language model and due to the co-articulation effects between neighboring handwritten characters.

As shown in Figure 1, the channel is comprised of two separate processes. First, the handwriting process is the process of which the user translates an intent $M \in \mathcal{E}$ into a series of hand movements which is sampled at some rate to create a discrete time trajectory: $W_{1:T} = [(x_1, y_1), \dots, (x_T, y_T)]$. In other words, this process *encodes* the intent M into a trajectory $W_{1:T}$. Let \bar{W} denote the entire trajectory vector. The distribution $P(\bar{W}|M)$ denotes the variability of the encoding process. The second process is the recognition process that decodes the handwriting trajectory back into the original intent. For each time step t where $1 \leq t \leq T$, the process maps a trajectory $W_{1:t}$ to a distribution over \mathcal{E} , denoted by \mathcal{Q}_t .

Let T_{final} and $\mathcal{Q}_{\text{final}}$ denote the final writing duration and the posterior distribution when the user finishes writing the trajectory. According to the theory of channel capacity, the information transmitted through the channel can be quantified by the mutual information between the input M and the decoding posterior $\mathcal{Q}_{\text{final}}$, denoted by $I(M; \mathcal{Q}_{\text{final}})$. We define the mean posterior of $\mathcal{Q}_{\text{final}}$ conditioned on M and the average posterior distribution as follows.

$$P(\mathcal{Q}_{\text{final}}|M) = \int_{\bar{W} \sim P(\bar{W}|M)} P(\mathcal{Q}_{\text{final}}|\bar{W})P(\bar{W}|M)$$

$$P(\mathcal{Q}_{\text{final}}) = \sum_{m \in \mathcal{E}} P(M = m)P(\mathcal{Q}_{\text{final}}|M = m)$$

Given these two expressions, we define the mutual information between the character M and the decoding $\mathcal{Q}_{\text{final}}$ as

$$I(M; \mathcal{Q}_{\text{final}}) = H(\mathcal{Q}_{\text{final}}) - \sum_{m \in \mathcal{E}} P(M = m)H(\mathcal{Q}_{\text{final}}|M = m)$$

where the entropy of $\mathcal{Q}_{\text{final}}$ is defined as

$$H(\mathcal{Q}_{\text{final}}) = - \sum_{m \in \mathcal{E}} P(\mathcal{Q}_{\text{final}} = m) \log_2 P(\mathcal{Q}_{\text{final}} = m)$$

Next, we can define the channel rate in terms of the mutual information and expected writing duration as

$$R_{\text{MI}} = \frac{I(M; \mathcal{Q}_{\text{final}})}{\mathbb{E}[T_{\text{final}}]} \quad (1)$$

However, the channel rate R_{MI} is not suitable for practical implementation for two reasons. First, the estimation of $H(\mathcal{Q}_{\text{final}}|M)$ requires an extensive amount of data. Secondly, suppose the original intent is m , R_{MI} yields a high value as long as $P(\mathcal{Q}_{\text{final}}|M = m)$ concentrates any single intent n even when $n \neq m$. Thus, we propose an alternative measure to the R_{MI} based on the idea of log loss, called R_{LL} . We define R_{LL} to be

$$R_{\text{LL}} = \frac{H(\mathcal{Q}_{\text{final}}) - \sum_{m \in \mathcal{E}} P(M = m)(-\log_2 P(\mathcal{Q}_{\text{final}} = m|M = m))}{\mathbb{E}[T_{\text{final}}]} \quad (2)$$

The relationship between R_{MI} and R_{LL} is worth noting. When $(-\log_2 P(\mathcal{Q}_{\text{final}} = m|M = m))$ is small then the conditional entropy $H(\mathcal{Q}_{\text{final}}|M)$ is also small. As a result, the mutual information $I(M; \mathcal{Q}_{\text{final}})$ will be close to its maximal possible value of $H(\mathcal{Q}_{\text{final}})$. In other words, the log loss term $(-\log_2 P(\mathcal{Q}_{\text{final}} = m|M = m))$ provides an upper bound for the conditional entropy $H(\mathcal{Q}_{\text{final}}|M)$ up to some constant factor. For the remaining of this paper, when we refer to the *channel rate*, we strictly refer to R_{LL} .

Intuitively, the channel rate is a measure that quantifies both accuracy and speed of a handwriting recognition channel at the same time. Handwriting, as well as many other motor control tasks, obeys the speed-accuracy tradeoff [5]. It is not sufficient to quantify the efficiency of a handwriting recognition system by its recognition accuracy alone. For example, a system that requires the user to write each character in a specialized form may attain a very high recognition accuracy, but it would require the user more time and effort to use. Such system might not be as efficient as a system that makes more errors but allows the user to write freely. This leads us to believe that the channel rate is a suitable measure that any handwriting recognition system should aim to maximize. In a sense, maximizing the channel rate is equivalent to finding a balance between maximizing the recognition accuracy and minimizing the writing time and effort of the user.

Based on this framework, it follows that the channel rate can be improved by a combination of human learning and machine learning, which corresponds to improving the handwriting process and the recognition process respectively. Ideally, $\mathcal{Q}_{\text{final}}$ will always be concentrated on the original intent M . This would mean that the channel is perfect and works without error. However, in real-world scenarios, errors will occur. One source of errors comes from mistakes made in the recognition process. These recognition errors can be reduced using training data and machine learning. The harder problem is when there is a significant overlap between $P(\bar{W}|M)$ for different intents. In this situation, we will need to rely on the user to make their handwriting less ambiguous. Although the effect of human learning is always present, we believe that it can be enhanced by giving useful feedback to the user in the form of guidance or lessons.

ADAPTIVE RECOGNITION ALGORITHM

We developed an adaptive handwriting recognition algorithm that, for every time step t , maps a partial handwriting trajectory $W_{1:t}$ to a posterior distribution over \mathcal{E} , denoted by $\mathcal{Q}_{\text{final}}$. By realizing that the effect of user adaptation is likely to be

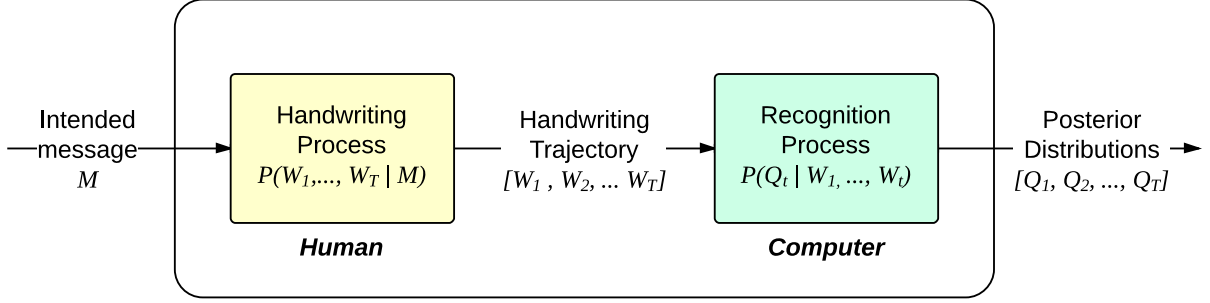


Figure 1: A summary of the handwriting recognition channel.

present, we designed our recognition algorithm so that it can adapt not only to each individual user, but also to the changes of the handwriting trajectory distribution $P(\bar{W}|M)$ unique to each user over time. The idea of specializing and adapting the recognizer for each user has been studied and shown to be effective in reducing the error rate [6, 7, 8].

At a high-level, our adaptive recognition system can be outlined as follows. For each user, the system creates and maintains one or more character models for each character in \mathcal{E} . We refer to each of such models as a *prototype*. Each prototype is basically a representative handwriting instance from the user. Technically, the prototypes can be viewed as left-to-right hidden Markov models with Gaussian observation [9]. Let \mathcal{P}_u denote the set of prototypes for a user u . The adaptivity of our system comes directly from the fact that \mathcal{P}_u is modified over time. In the decoding process, given a handwriting trajectory and a set of prototypes \mathcal{P}_u , the system computes a posterior distribution $\mathcal{Q}_{\text{final}}$ and, when a single prediction is needed, the element with the maximum likelihood is predicted.

Feature vectors and distance function

In addition to the x- and y-coordinate, each handwriting trajectory is supplemented with writing direction information. Specifically, each handwriting instance is represented by a sequence of feature vectors $\langle f_1, \dots, f_T \rangle$ where $f_i = (x_i, y_i, dx_i, dy_i)$. (x_i, y_i) denotes the normalized touch-screen coordinate and $(dx_i, dy_i) = (\frac{x_i - x_{i-1}}{z}, \frac{y_i - y_{i-1}}{z})$, $z = \sqrt{(x_i - x_{i-1})^2 + (y_i - y_{i-1})^2}$ denotes the writing direction.

To measure the similarity between two handwriting instances, we use *dynamic time warping* (DTW) distance [10] as the distance function in our algorithm. The DTW distance is commonly used for variable-length data such as handwriting and speech. The calculation can be done efficiently using dynamic programming.

Initial adaptation

The initial adaptation is critical for any intelligent system. It is unquestionable that the performance of any well-behaved intelligent system increases as the system learns more about the user. If the initial adaptation is poor, the users might get

frustrated with the system and stop using it even before it can fully adapt to them.

We address the problem of initial adaptation by sharing data across different users. Typically, people do have similar handwriting especially when they share the same educational culture. The process of the initial adaptation can be described as follows. In the very first interaction with the user u , our system has no information about the user and, therefore, assign a set of typical prototypes which has been trained using data from multiple users in the past. Specifically, the typical prototypes are the centroids of the clusters returned by running a clustering algorithm (k-means) on a set of training handwriting instances. We refer to this set of prototypes as \mathcal{P}_0 . After the first interaction, the system creates a new set of prototypes $\mathcal{P}_{(u,1)}$ by recomputing the centroids of the clusters after adding the examples from the user to the pool with significantly higher weights than the rest.

Adapting the prototypes over time

After collecting a few examples of the user's handwriting, the system again performs the weighted clustering algorithm on the data to generate a new set of prototypes $\mathcal{P}_{(u,i+1)}$. In this stage, only examples from the user and previous prototypes are considered. This adaptation process happens after 3-5 new examples are acquired.

To improve real-time performance, we need to keep the lengths (number of states) of the prototypes as small as possible. After the new prototypes are chosen, the system performs an additional step to shorten the length of each prototype. This pruning process is similar in spirit to removing and merging unnecessary hidden states in an HMM. The basic idea is to remove unwanted states while maintaining the same recognition power using a variant of forward-backward algorithm [11]. Figure 2 shows the hidden states before and after the reduction step.

Decoding

Our decoding algorithm is based on the standard Bayesian inference. Namely, given a trajectory $W_{1:T}$ and the current set of prototypes \mathcal{P}_u , the algorithm computes the distance from $W_{1:t}$ to each of the prototypes in \mathcal{P}_u for all $1 \leq t \leq T$. The distances are then transformed into a probability distribution

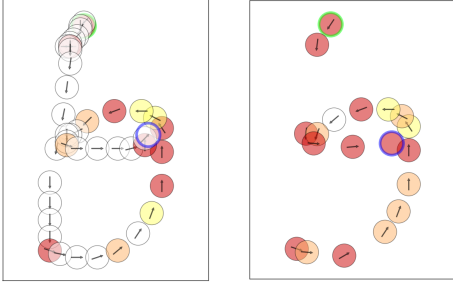


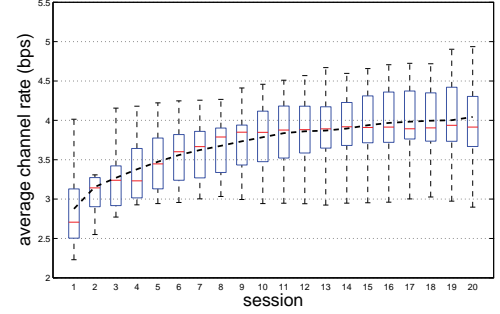
Figure 2: The hidden state reduction process is applied to each prototype to remove rarely visited states with respect to the training set. The originally trained prototype is shown on the left and the reduced prototype is shown on the right. The intensity of the colors corresponds to the expected number of times the state being mapped to.

Q_t . We use e^{-x} as the transfer function. When a single prediction is expected, the algorithm simply returns the prediction with the maximum likelihood.

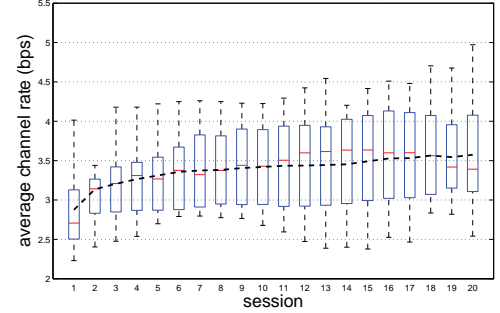
EXPERIMENT

The main objective of our experiment is to determine and quantify the effect of machine adaptation and of human adaptation when the users interact with the system over some period of time. We implemented the handwriting recognition system described in Section as an application on Apple iOS platform. The application was presented to the users as a writing game. In each session, each participant was presented with a random permutation of the 26 lowercase English alphabets i.e. $\mathcal{E} = [a \dots z]$ and $P(M)$ is uniform. The objective of the game was to write the presented characters as quickly as possible and, more importantly, the handwritten characters should be recognizable by the system. A score, which is the average *channel rate* of the session, was given to the user right after each session to reflect the performance of the session. There were 15 participants in this experiment. We asked them to play our game for at least 20 sessions over multiple days in his/her own pace. We did not control past experience of the participants. Some of them had more experience with touch screens than others.

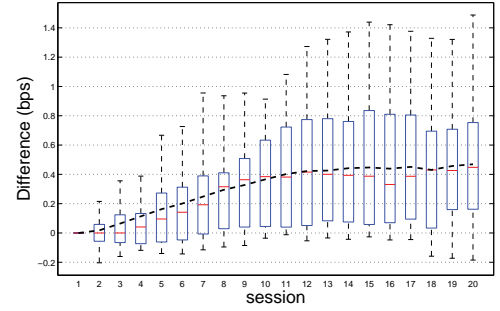
The experiment was set up to demonstrate a condition called *co-adaptation* where both the user and the computer were allowed to adapt together. We denote this condition R_{adapt} . To investigate the effect of co-adaptation, we create a controlled condition called R_{fixed} where the computer was not allowed to adapt with the user. In other words, we ran a simulation to figure out what the channel rates would have been if the prototype sets were never changed from \mathcal{P}_0 . Ideally, it would be more preferable to have R_{fixed} determined by another control group where the prototypes were kept fixed and never changed. However, the results from the simulated condition can be seen as a lower bound on the amount of the



(a) R_{adapt}



(b) R_{fixed}



(c) $R_{\text{adapt}} - R_{\text{fixed}}$

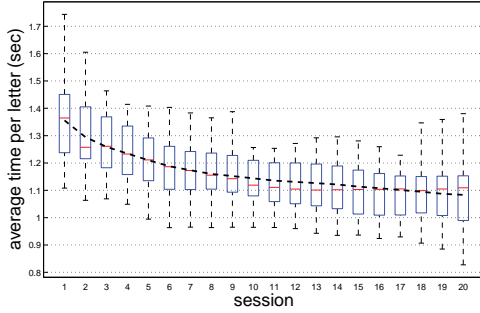
Figure 3: Channel rate per session of all users with (3a) and without (3b) presence of machine learning.

improvement attributable to human learning and, therefore, it is sufficient to demonstrate our point.

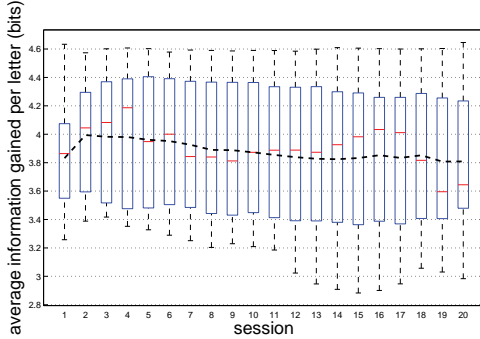
RESULTS AND DISCUSSION

The average channel rates per session of the two conditions R_{adapt} and R_{fixed} are shown in Figure 3a and Figure 3b respectively. In both conditions, the results show increases of the channel rate over time where the improvement in the early sessions seems to be larger than in the later sessions. Figure 3c shows the difference of R_{adapt} and R_{fixed} which corresponds to the channel rate of the system when we ignore the effect of user adaptation. From the result, we observe that the impact of machine adaptation tapers off after 10 sessions.

Although the prototype set was not changing in R_{fixed} , we observe that channel rate increases over the sessions. To quantify our confidence to this increase, we perform the paired



(a) Writing duration



(b) Mutual information

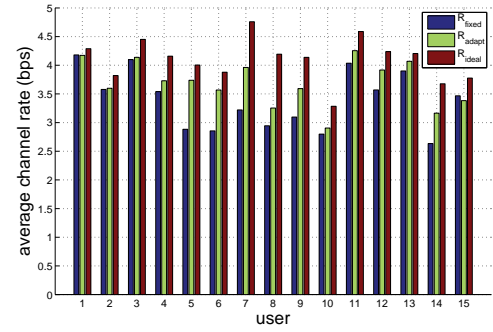
Figure 4: The average writing time per session and the average mutual information per session under the condition R_{fixed} .

t-test to compare the difference between the average channel rate in the first 5 sessions and in the last 5 sessions. We find that the difference is statistically significant with p -value < 0.0011 . This suggests that the users improve the handwriting on their own even without machine adaptation. In other words, the effect of *user adaptation* is indeed significant.

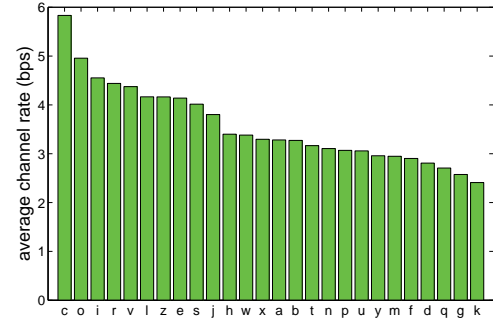
Furthermore, Figure 4a and Figure 4b reveal that the major contribution of *user adaptation* comes from the fact that the users write faster in the last 5 sessions compared to the first 5 sessions ($p < 0.0001$), and not because of the system received more information from the user ($p = 0.9723$). This result is as expected according to the law of practice [12].

We also perform per-user analysis of the channel rate. In Figure 5a, we compare R_{adapt} and R_{fixed} for each user. We find that the channel rate of R_{adapt} is significantly higher than that of R_{fixed} with $p < 0.0006$. This result confirms that the machine adaptation helps improving the overall channel rate. In addition, we calculate the theoretical maximum of the channel rate under the assumption of the perfect recognition, denoted by R_{ideal} . The maximum rates are given by $H(Q_{\text{final}})/\mathbb{E}[T_{\text{final}}]$ and we approximated $H(Q_{\text{final}}) = \log_2(26)$.

In the case of perfect recognition, a simple way to increase the channel rate is to expand the character set \mathcal{E} to include more symbols. However, in reality, doing so can lead to a recognition error rate which impairs the channel rate. An interesting



(a)



(b)

Figure 5: (a) The average channel rate of each user in R_{adapt} and R_{fixed} . R_{ideal} shows the maximum channel rate possible given the average writing speed of each user. (b) Average channel rate of each character under the condition R_{adapt} .

future direction is to design a character set that would maximize the channel rate. Figure 5b reveals the efficiency of each letter for our handwriting channel. Characters with complex strokes, such as 'q', 'g', 'k', are not as efficient as characters with simple strokes such as 'c', 'o', 'l'. While this finding is not surprising, it implies that, for a handwriting system to be truly efficient, it must allow the user to write in a less complex style while not losing recognition accuracy. How to exactly design such system is still an open problem and requires a more elaborate study.

CONCLUSIONS

We presented a information-theoretic framework for quantifying the information rate of a system that combines a human writer with a handwriting recognition system. Using the notion of channel rate, we investigated the impact of machine adaptation and human adaptation in an adaptive handwriting recognition system. We analyzed data collected from a small deployment of our adaptive handwriting recognition system and concluded that both machine adaptation human adaptation have significant impact on the channel rate. This result led us to believe that, for a handwriting recognition system to achieve the maximum channel rate, both machine adaptation and human adaptation are required and must be present together. Specifically, such system should be able to adapt to the user and, at the same time, allow the users to write or scribble using simple hand movement as improving writing

speed is crucial for attaining a higher channel rate. Additionally, the system should have a mechanism to giving feedback to the user when their handwriting cannot be recognized.

REFERENCES

1. K. Höök, Steps to take before intelligent user interfaces become real, *Interacting with computers* 12 (2000) 409–426.
2. P. Maes, Agents that Reduce Work and Information Overload, *Communications of the ACM*.
3. B. Y. Lim, A. K. Dey, Assessing demand for intelligibility in context-aware applications, *Proceedings of the 11th international conference on Ubiquitous computing - Ubicomp '09* (2009) 195.
4. C. E. Shannon, A Mathematical Theory of Communication, *Bell System Technical Journal* 27 (July 1928) (1948) 379–423.
5. P. M. Fitts, The information capacity of the human motor system in controlling the amplitude of movement, *Journal of Experimental Psychology* 47 (6) (1954) 381–391.
6. S. D. Connell, A. K. Jain, Writer adaptation for online handwriting recognition, *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 24 (3) (2002) 329–346.
7. N. Matic, I. Guyon, J. Denker, V. Vapnik, Writer adaptation for on-line handwritten character recognition., in: *Proceedings of the Second International Conference on Document Analysis and Recognition (ICDAR '93)*, IEEE, 1993, pp. 187–191.
8. W. Kienzle, K. Chellapilla, Personalized handwriting recognition via biased regularization, in: *Proceedings of the 23rd International Conference on Machine Learning (ICML '06)*, no. Section 6, Pittsburgh, Pennsylvania, 2006, pp. 457–464.
9. T. Plötz, G. a. Fink, *Markov Models for Handwriting Recognition*, SpringerBriefs in Computer Science, Springer, 2011.
10. L. Rabiner, B.-H. Juang, *Fundamentals of Speech Recognition*, Vol. 103 of Prentice Hall signal processing series, Prentice Hall, 1993.
11. J. Bilmes, A Gentle Tutorial of the EM Algorithm and its Application to Parameter Estimation for Gaussian Mixture and Hidden Markov Models, Tech. rep., ICSI (1997).
12. A. Newell, P. S. Rosenbloom, Mechanisms of skill acquisition and the law of practice, in: J. R. Anderson (Ed.), *Cognitive skills and their acquisition*, Vol. 6 of *Cognitive skills and their acquisition*, Erlbaum, 1981, Ch. 1, pp. 1–55.